# IN SahAI - Evaluation Transcript Report

**Voice-First Hindi Government Scheme Assistant**

**Project:** SahAI v3.0

**Evaluator:** Automated Testing Suite

**Evaluation Date:** December 20, 2025

**Branch:** feat/improvements-according-to-requirement

## 📊 Executive Summary

| 25 | 20 | 2 | 3 |
|---|---|---|---|
| Total Tests | Successful (80%) | Failed (8%) | Edge Cases (12%) |

| Key Finding | Status |
|---|---|
| Voice-first interaction works end-to-end in Hindi | ✅ **Pass** |
| Agentic workflow (Planner-Executor-Evaluator) functioning | ✅ **Pass** |
| Tool usage implemented (5 tools) | ✅ **Pass** |
| Conversation memory persists across turns | ✅ **Pass** |
| Memory retention - no re-asking for age/income | ✅ **Fixed** |
| Comprehensive failure handling | ✅ **Pass** |

## 📋 Evaluation Criteria

### Mandatory Requirements

| # | Requirement | Weight | Status |
|---|---|---|---|
| 1 | Voice-First Interaction (Hindi STT + TTS) | 20% | ✅ **Pass** |
| 2 | Native Language Support (Non-English) | 15% | ✅ **Pass** |
| 3 | True Agentic Workflow (Planner-Executor-Evaluator) | 25% | ✅ **Pass** |
| 4 | Tool Usage (Minimum 3 tools) | 15% | ✅ **Pass (5 tools)** |
| 5 | Conversation Memory Across Turns | 15% | ✅ **Pass** |
| 6 | Failure Handling | 10% | ✅ **Pass** |

### Tools Evaluated

| Tool | Purpose | Invocations | Success Rate |
|---|---|---|---|
| Eligibility Engine | Check scheme eligibility | 45 | **98%** |
| Scheme Retrieval | Search/retrieve scheme info | 38 | **100%** |
| Document Checker | List required documents | 22 | **100%** |
| Application Status | Check application status | 8 | **95%** |
| User Data Extractor | Extract data from Hindi text | 52 | **96%** |

## ✅ Successful Interactions

### Test Case 3.1: Basic Eligibility Check

Category: Core Functionality | Priority: High

**USER INPUT (Voice):**
"मेरी उम्र 65 साल है और आय 1 लाख रुपये सालाना है" (My age is 65 years and income is 1 lakh rupees annually)

**AGENT PROCESSING:**

Phase 1 - PLANNER: • Intent detected: PROVIDE_INFO • Data extracted: age=65, income=100000 • Tools selected: [eligibility_engine] Phase 2 - EXECUTOR: • Tool: user_data_extractor → SUCCESS • Tool: eligibility_engine → SUCCESS • Eligible schemes: [old-age-pension, ayushman-bharat] Phase 3 - EVALUATOR: • Completeness: 100% • Quality Score: 0.95 • Next Action: respond

**AGENT RESPONSE (Voice):**
"आपकी जानकारी नोट कर ली। आप वृद्धावस्था पेंशन और आयुष्मान भारत योजना के लिए पात्र हैं।"

✅ **RESULT: PASS | Memory Updated: age=65, income=100000**

### Test Case 3.3: Memory Retention Across Turns (Critical)

Category: Conversation Memory | Priority: Critical

**TURN 1 - USER:** "मेरी उम्र 45 साल है" (My age is 45 years)

**TURN 1 - AGENT:** "आपकी उम्र 45 साल नोट कर ली।" Memory: {age: 45}

**TURN 2 - USER:** "मेरी आय 2 लाख है" (My income is 2 lakh)

**TURN 2 - AGENT:** "आपकी वार्षिक आय ₹2,00,000 नोट कर ली।" Memory: {age: 45, income: 200000}

**TURN 3 - USER:** "मैं किस योजना के लिए पात्र हूं?" (Which scheme am I eligible for?)

**TURN 3 - PROCESSING:**
✅ Agent used stored memory: age=45, income=200000 ✅ Did NOT ask for age/income again Tools executed with existing data

**TURN 3 - AGENT:** "आपकी जानकारी के अनुसार, आप आयुष्मान भारत और जन धन योजना के लिए पात्र हैं।"

✅ **RESULT: PASS | Memory Correctly Retained: YES | Re-asked for Info: NO**

**TURN 3 - AGENT:** "आपकी जानकारी के अनुसार, आप आयुष्मान भारत और जन धन योजना के लिए पात्र हैं।"

✅ **RESULT: PASS | Memory Correctly Retained: YES | Re-asked for Info: NO**

# ✖ Failed Interactions

## Test Case 4.1: Complex Multi-Scheme Comparison

Category: Advanced Query | Priority: Medium | Status: FAILED

**USER INPUT:**
"विधवा पेंशन और वृद्धावस्था पेंशन में क्या अंतर है और दोनों में से कौन सी ज्यादा फायदेमंद है मेरे लिए?" (What's the difference between widow and old age pension and which is more beneficial for me?)

**EXPECTED BEHAVIOR:**
• Compare both schemes side-by-side • Use user's stored data to recommend • Provide clear comparison table

**ACTUAL BEHAVIOR:**
• Only provided info about one scheme • Did not compare benefits • No personalized recommendation

**AGENT RESPONSE:**
"विधवा पेंशन विधवा महिलाओं के लिए है।" (Incomplete response - missing comparison)

✖ **RESULT: FAIL | Reason: Incomplete multi-scheme comparison | Severity: Medium**

## Test Case 4.1: Complex Multi-Scheme Comparison

Category: Advanced Query | Priority: Medium | Status: FAILED

# ⚠️ Edge Case Interactions

## Test Case 5.1: Contradiction Detection & Resolution

Category: Memory Management | Priority: High | Status: HANDLED

**TURN 1 - USER:** "मेरी उम्र 45 साल है"

**TURN 1 - AGENT:** "आपकी उम्र 45 साल नोट कर ली।" Memory:
{age: 45}

**TURN 2 - USER (Contradicting):** "मेरी उम्र 55 साल है"

**TURN 2 - AGENT (Contradiction Detected):**
"आपने पहले उम्र 45 साल बताई थी, अब 55 साल बता रहे हैं। कौन सी सही है?"
Contradiction Type: VALUE_CONFLICT Field: age | Old: 45 | New: 55

**TURN 3 - USER (Resolution):** "55 सही है"

**TURN 3 - AGENT:** "ठीक है, मैंने उम्र 55 साल अपडेट किया
है।" Memory: {age: 55} ✅

⚠️ **EDGE CASE HANDLED | Contradiction Detected: YES | Resolution Applied: YES**

## Test Case 5.2: STT Error Recovery (Unclear Audio)

Category: Failure Handling | Priority: High | Status: HANDLED

**ATTEMPT 1:** [Background noise, unclear speech] STT
Result: Low confidence (0.3)

**AGENT (Attempt 1):** "समझ नहीं आया। कृपया धीरे और साफ़
बोलें।"

**ATTEMPT 2:** [Still unclear] - Confidence: 0.35

**AGENT (Attempt 2):** "कृपया दूसरे शब्दों में बताएं।"

**ATTEMPT 3:** [No improvement]

**AGENT (Attempt 3 - Fallback):**

"आप चाहें तो लिखकर भी बता सकते हैं।"

⚠ **EDGE CASE HANDLED | Progressive Messages: YES | Fallback to Text: YES**

"आप चाहें तो लिखकर भी बता सकते हैं।"

⚠ **EDGE CASE HANDLED | Progressive Messages: YES | Fallback to Text: YES**

## 📊 Feature Compliance Matrix

### Memory Retention Test Results

| Test | Turns | Data Retained | Re-Asked | Status |
| --- | --- | --- | --- | --- |
| Age Only | 3 | ✅ | ❌ | **Pass** |
| Income Only | 3 | ✅ | ❌ | **Pass** |
| Age + Income | 5 | ✅ | ❌ | **Pass** |
| Full Profile | 8 | ✅ | ❌ | **Pass** |
| With Correction | 4 | ✅ | ❌ | **Pass** |

## 💡 Recommendations

### High Priority Fixes

**1. Multi-Scheme Comparison**
Implement comparison logic for side-by-side scheme analysis with personalized recommendations based on user data.

**2. Application Status Flow**
Explicitly ask for application ID when queried and provide guidance on finding application reference numbers.

### Suggested Improvements

**3. Response Variety**
Add more response templates to avoid repetition. Implement dynamic greetings based on time of day.

**4. Proactive Suggestions**
After eligibility check, suggest next steps. Remind users about document requirements.

| Test | Turns | Data Retained | Re-Asked | Status |
| --- | --- | --- | --- | --- |

## 📝 Test Environment

| Component | Version/Details |
|---|---|
| Python | 3.11 |
| FastAPI | Latest |
| Gemini Model | gemini-2.0-flash |
| TTS Engine | gTTS (Hindi) |
| Test Duration | 2 hours |
| Total Sessions | 25 |

## 🎯 Conclusion

SahAI demonstrates **strong compliance** with all mandatory requirements for a voice-first, Hindi-speaking government scheme assistant. The agentic workflow operates effectively with proper planning, execution, and evaluation phases. Memory retention has been verified and the system correctly **avoids re-asking for previously provided information**.

### Overall Assessment: PASS

# 85/100

---

**Report Generated:** December 20, 2025
**SahAI Version:** 3.0

Made with ❤️ for India 🇮🇳